

Learning Efficient Feature Representation for Temporal Action Localization

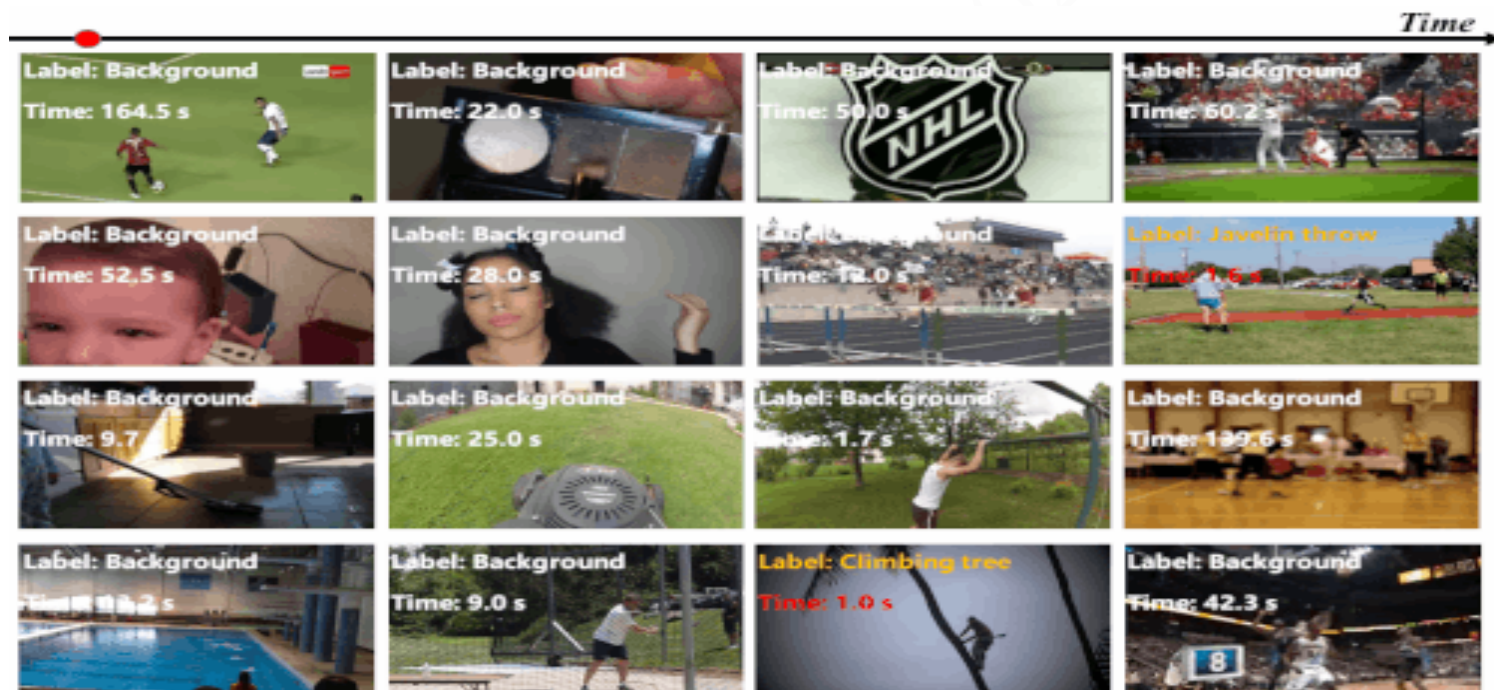
Team: Dahua_001

Chenglu Wu*, Xuefeng Yang*, Fuzhi Duan, Yanxun Yu,
Yayun Wang, Jun Yin

Zhejiang Dahua Technology Co., Ltd.

I. Task Description

Temporal action localization (TAL) requires to precisely locate the temporal boundaries of action instances and accurately classify the action instances into specific categories.



II. Evaluation Criteria

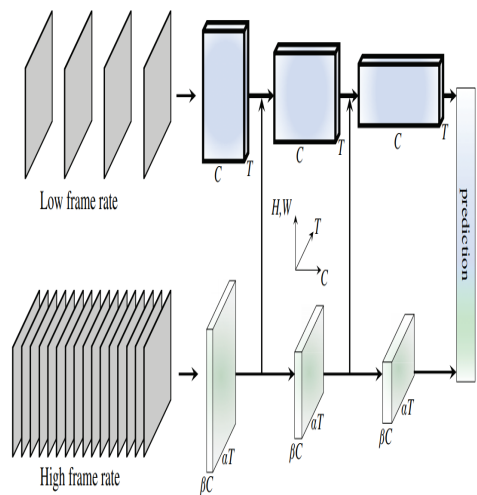
Interpolated Average Precision (AP) is used as the metric for evaluating the results on each activity category. Then, the AP is averaged over all the activity categories (mAP).

The official metric used in competition is the **average mAP**, which is defined as the mean of all mAP computed with tIoU thresholds between 0.5 and 0.95 with a step of 0.05.

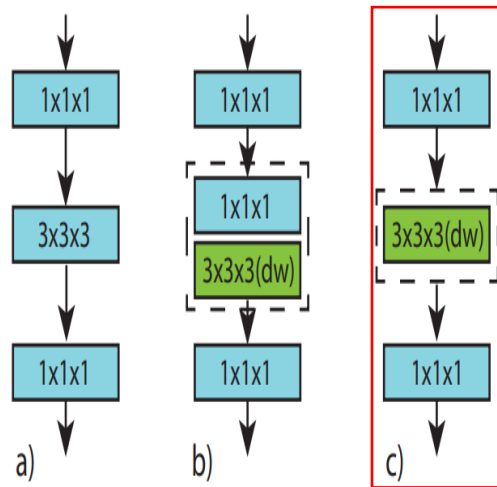
In addition to evaluate proposal quality, Average Recall (AR) under multiple tIoU thresholds are calculated. $AR@AN$ is defined as the AR under different average proposal numbers (AN), and we calculate the area under the AR vs. AN curve (**AUC**) as a metric of proposal.

III. Mainstream Algorithms

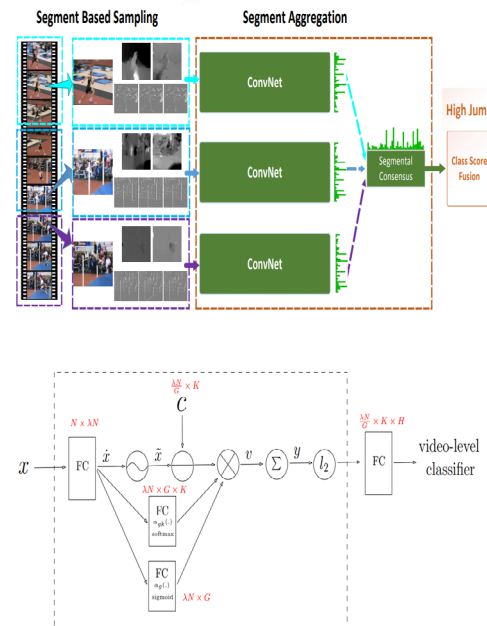
Video Classification.



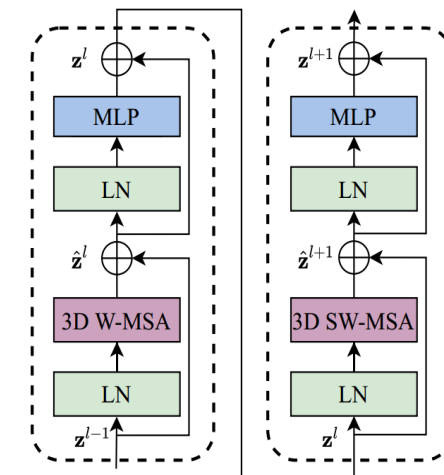
SlowFast



CSN



TSN&NeXtVLAD

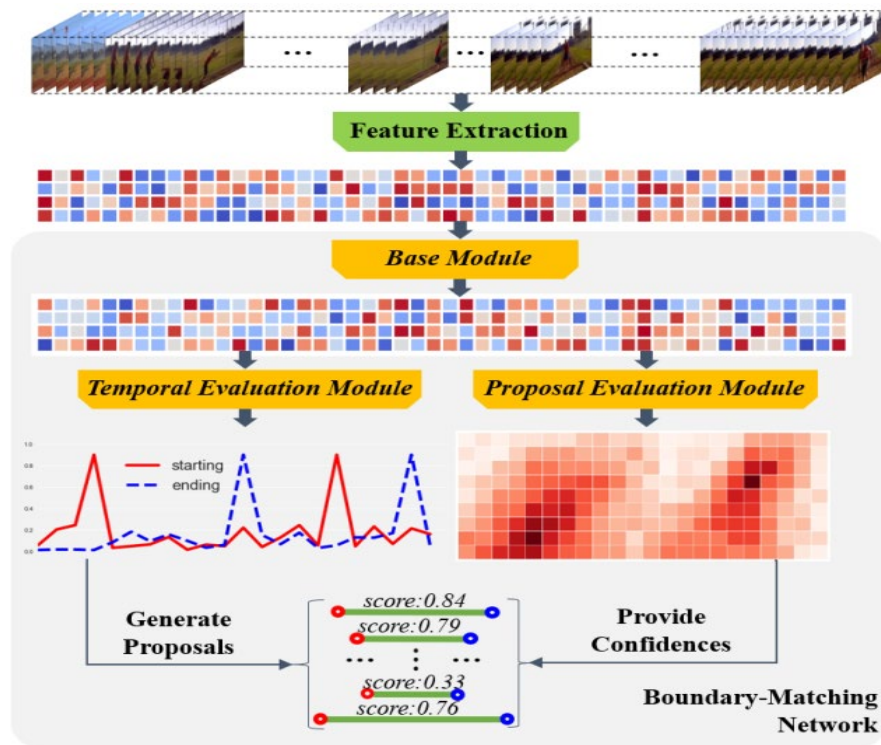


Video-swin-transformer

- [2]. Christoph Feichtenhofer, etc. Slowfast networks for video recognition. ICCV2019
- [3]. Du Tran, etc. Video classification with channel-separated convolutional networks. ICCV2019
- [4]. Limin Wang, etc. Temporal segment networks: Towards good practices for deep action recognition. ECCV2016
- [5]. Rongcheng Lin, etc. Nextvlad: An efficient neural network to aggregate frame-level features for large-scale video classification. ECCV-Workshops 2018
- [6]. Ze Liu, etc. Video swin transformer. arXiv preprint arXiv:2106.13230, 2021.

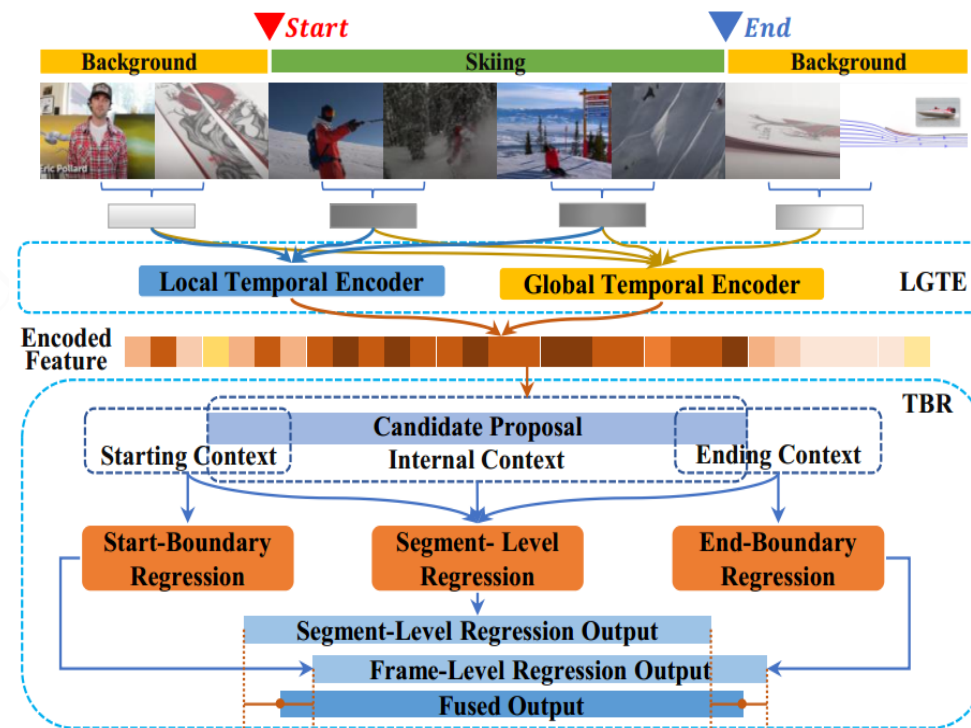
III. Mainstream Algorithms

Proposal Generation.



BMN

Proposal Refinement.



TCAnet

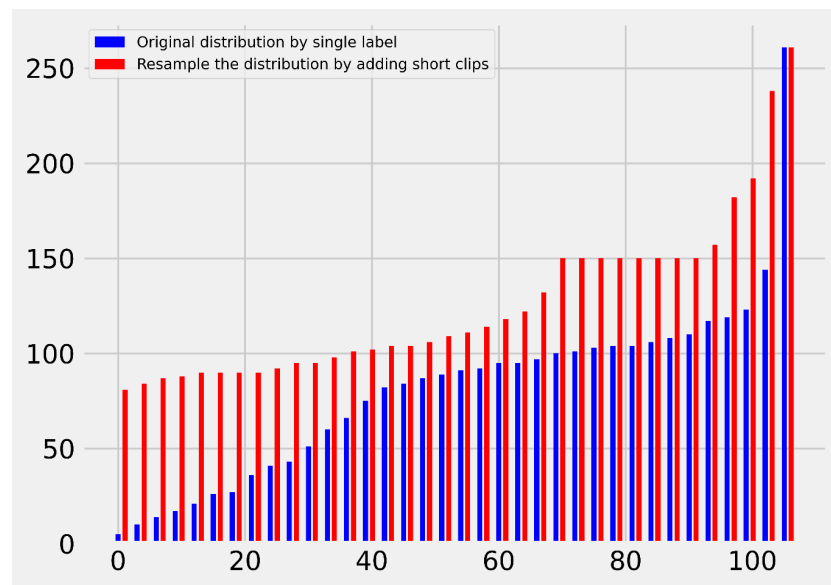
[7]. Tianwei Lin, etc. Bmn: Boundary-matching network for temporal action proposal generation. ICCV2019

[8]. Zhiwu Qing, etc. Temporal context aggregation network for temporal action proposal refinement. CVPR2021

IV. Our Approach

Dataset cleansing.

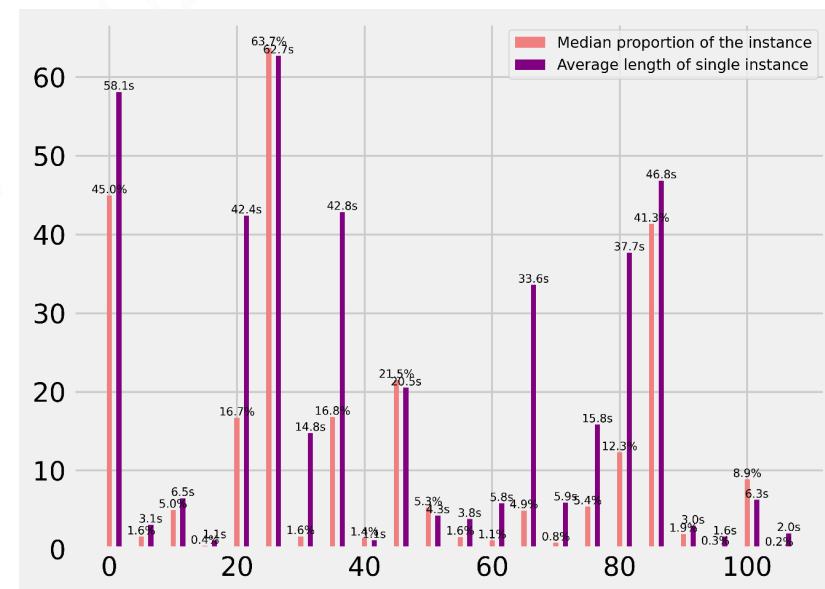
- Class Skew Balance



Original distribution VS. resample distribution.

Blue bar is the distribution of original number of different categories.
Red bar is the number of unbalanced video clips by resampling.

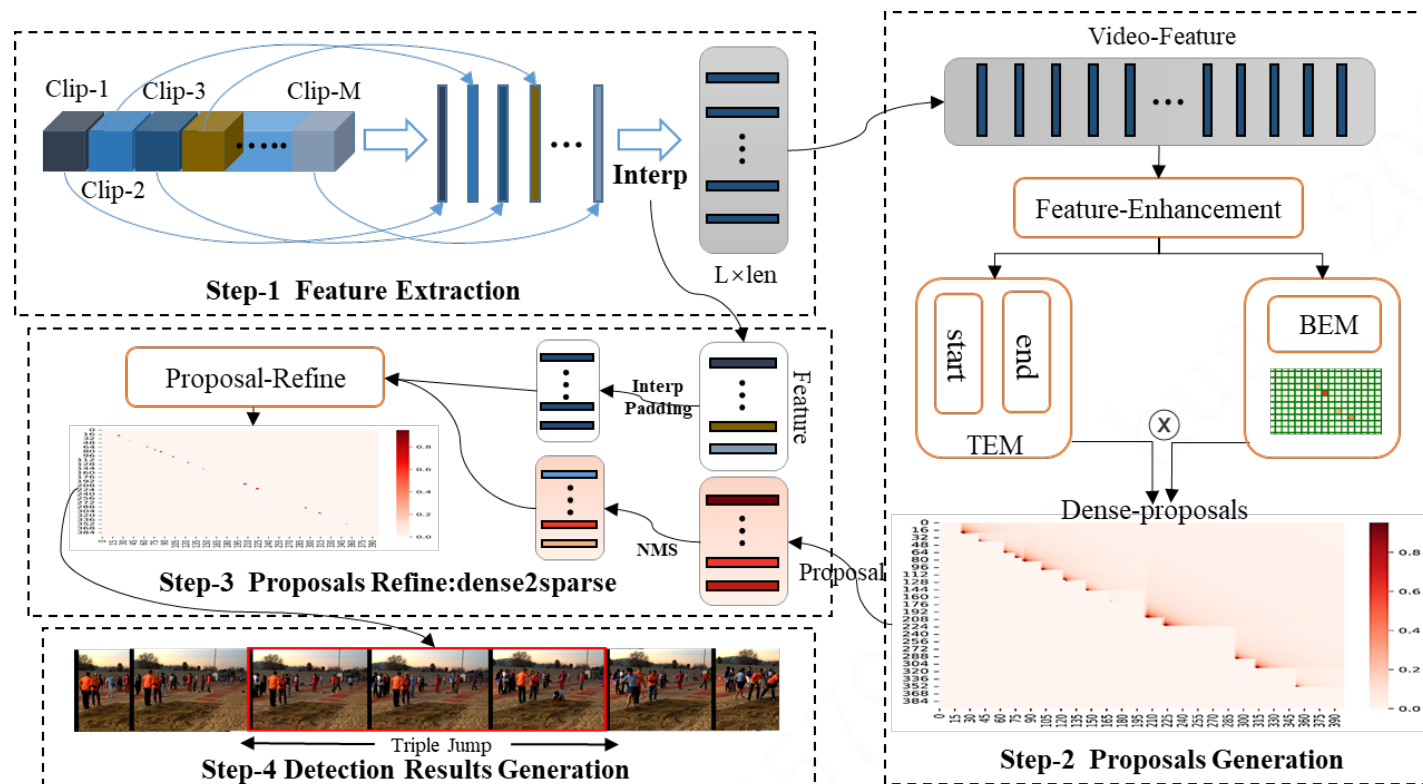
- Outlier Identification & Cleansing



Median proportion and the average length of instances.

Light coral bar is the median percentage of the total video duration for a single instance per category. Purple bar is the average length of single instance of each category.

IV. Our Approach



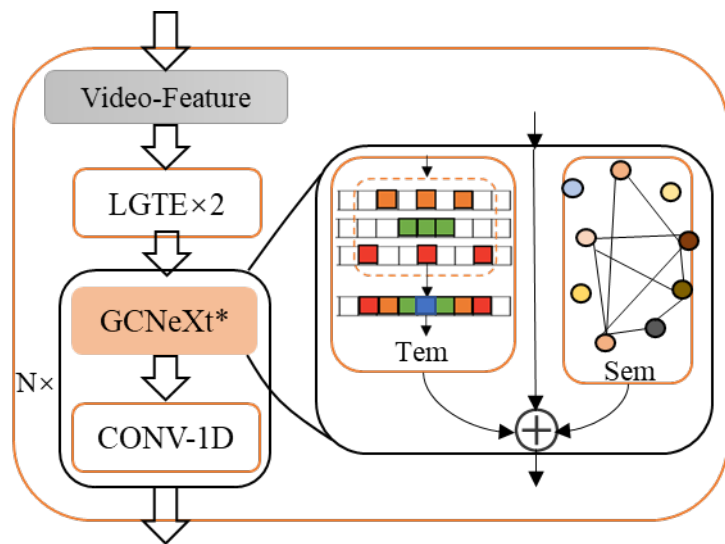
Pipeline:

- Feature Extraction.
- Temporal Proposal
- Proposal Refine:dense2sparse.
- Detection Results Generation.

Learning Efficient Features Representation for TAL (overall)

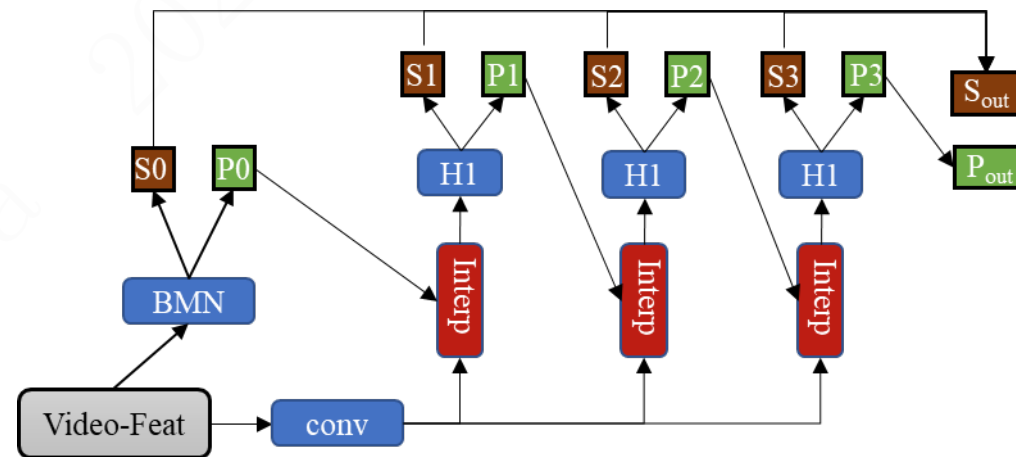
IV. Our Approach

Feature Enhancement module:



The improved base-feature layer.

- Encode local and global temporal relationships.
- Deform the receptive field and enhance the aggregation of context information



The architecture of cascade proposals refinement.

- Pick samples corresponding to a specific iou for training.
- Boost the performance of the proposal prediction gradually.

V. Experiments Results

Classification Results:

Method	SlowFast	CSN	TSN	NeXtVLAD	Video-swin-transformer
backbone	Res101+50	Res152	Swin-Base	Swin-Base	Swin3D-Base
head	SlowFastHead	I3DHead	TSNHead	NextVLADHead	I3DHead
clip_len	32	32	1	2	32
frame_interval	2	2	1	1	3
num_clips	1	1	32	32	1

Training details for video classification networks.

Model	Top-1	Top-2	Top-3	Top-5
TSN	81.54	91.46	94.19	96.87
SlowFast (nc=8)	84.92	94.11	96.47	98.63
CSN (nc=8)	85.70	94.39	96.96	98.78
Video-SwinB (nc=5)	87.61	94.90	97.41	98.82
Video-SwinB (nc=8)	87.90	95.23	97.84	99.04
NeXtVLAD	87.21	94.87	97.24	98.70
Ensemble	90.03	96.72	98.43	99.41

We used **NeXtVLAD**, **CSN** and **Video-SwinB** to ensemble the model.

V. Experiments Results

Proposal Results:

Video feat	L	AR@1	AR@5	AR@10	AR@100	AUC
I3D	100	4.92	10.35	13.38	24.64	19.57
	200	4.87	10.40	13.75	27.78	21.28
TSN-K700	200	5.15	11.44	15.10	29.61	23.07
	250	5.05	10.90	14.41	28.64	22.04
Slowonly-k700	200	5.09	11.19	14.86	29.25	22.67
TSN-full-2048	200	5.42	12.15	16.08	31.29	24.53
SwinB-k600	200	5.80	12.98	16.94	31.73	25.05
	256	5.82	12.90	17.09	32.26	25.32
	325	6.21	14.23	18.75	35.00	27.69
SwinB-full-1024	256	6.07	13.97	18.60	33.96	27.28
	325	6.21	14.44	19.32	36.84	29.25
	400	6.49	15.13	20.30	38.80	30.81
	450	6.34	15.20	20.32	38.61	30.67

- The AUC of the proposal at L400 is higher than L325 and L450 from the validation results of different temporal scales.
- The feature extracted by the SwinB model showed a significant improvement compared to the other features.

The L denotes the length of video feature.

V. Experiments Results

Detection Results:

BMN	LGTE	GCNeXt	Dilate	TCAnet	NMS	Cascade	Average-mAP(%)	Promotion
✓							17.32	-
✓	✓						17.61	+1.67%
✓	✓	✓					18.21	+3.41%
✓	✓	✓	✓				18.63	+2.31%
✓	✓	✓	✓	✓			19.14	+2.74%
✓	✓	✓	✓	✓	✓		21.19	+10.71%
✓	✓	✓	✓	✓	✓	✓	22.05	+4.06%

Influence of different modules on the performance of FineAction.

- The BMN network trained with the standardized video feature length of 400 was the baseline.
- The detection result achieves 22.05% on the validation set and 23.35% on the test set in terms of average mAP.

VI. Conclusion

- We find that the model using the video clips for action recognition has a greater performance on proposals than the model using single frame.
- Increasing the grid size of BMN can further improve the recognition accuracy of short actions in the FineAction. However, with the increase of grid size, the total number of parameters of the model will increase exponentially, and the model convergence epoch will move backward as well.
- Considering the cumbersome nature of this method, we will improve the one-stage TAL framework to locate and identify the extremely short instances in the future.

THANKS